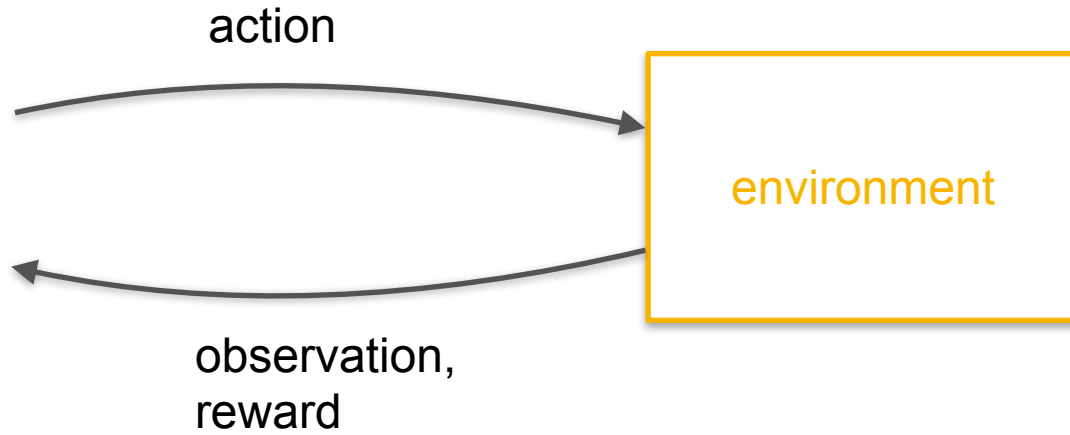# Framework Design for Deep Reinforcement Learning

先端人工知能論２／第８回／演習

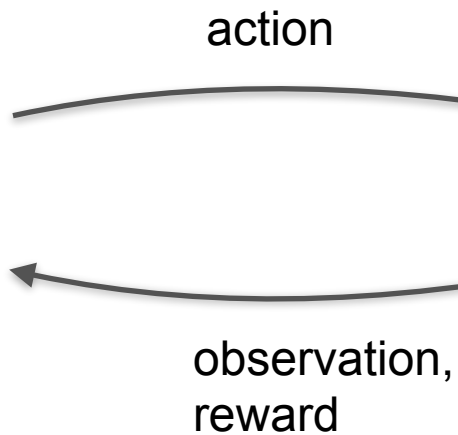2017.11.21 (Tue)

片岡 俊基 (Preferred Networks)

# Recall

# Recall

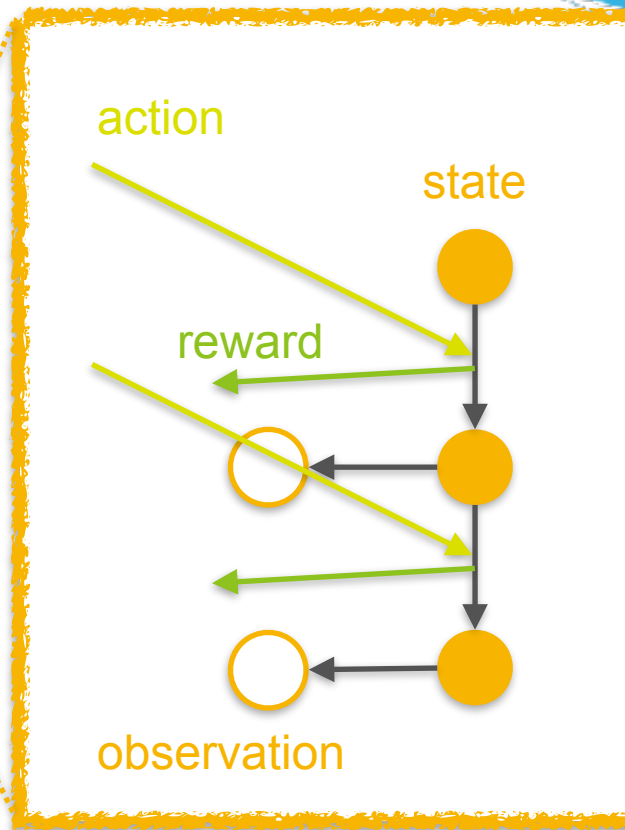# Case Study: Deep Q-Network

- Algorithm in the DQN paper [Mnih+ 2013]

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$
        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
    **end for**
**end for**

Second, learning directly from consecutive samples is inefficient, due to the strong correlations between the samples; randomizing the samples breaks these correlations and therefore reduces the variance of the updates. Third, when learning on-policy the current parameters determine the next data sample that the parameters are trained on. For example, if the maximizing action is to move left then the training samples will be dominated by samples from the left-hand side; if the maximizing action then switches to the right then the training distribution will also switch. It is easy to see how unwanted feedback loops may arise and the parameters could get stuck in a poor local minimum, or even diverge catastrophically [25]. By using experience replay the behavior distribution is averaged over many of its previous states, smoothing out learning and avoiding oscillations or divergence in the parameters. Note that when learning by experience replay, it is necessary to learn off-policy (because our current parameters are different to those used to generate the sample), which motivates the choice of Q-learning.

In practice, our algorithm only stores the last $N$ experience tuples in the replay memory, and samples uniformly at random from $\mathcal{D}$ when performing updates. This approach is in some respects limited since the memory buffer does not differentiate important transitions and always overwrites with recent transitions due to the finite memory size $N$. Similarly, the uniform sampling gives equal importance to all transitions in the replay memory. A more sophisticated sampling strategy might emphasize transitions from which we can learn the most, similar to prioritized sweeping [17].

**4.1 Preprocessing and Model Architecture**

Working directly with raw Atari frames, which are $210 \times 160$ pixel images with a 128 color palette,

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$
        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
    **end for**
**end for**

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$

Initialize action-value function $Q$ with random weights

**for** episode $= 1, M$ **do**

    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

    **for** $t = 1, T$ **do**

        With probability $\epsilon$ select a random action $a_t$

        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$

        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$

        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

    **end for**

**end for**

**Algorithm 1** Deep Q-learning with Experience Replay

obs. $\mapsto$ state

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$
        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
    **end for**
**end for**

**Algorithm 1** Deep Q-learning with Experience Replay

   Initialize replay memory $\mathcal{D}$ to capacity $N$
   Initialize action-value function $Q$ with random weights
   **for** episode $= 1, M$ **do**
      Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
      **for** $t = 1, T$ **do**
         With probability $\epsilon$ select a random action $a_t$
         otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
         Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
         Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
         Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$
         Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$
         Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
         Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
      **end for**
   **end for**

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$

Initialize action-value function $Q$ with random weights

**for** episode $= 1, M$ **do**

    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

    **for** $t = 1, T$ **do**

        With probability $\epsilon$ select a random action $a_t$

        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$

        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$

        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

    **end for**

**end for**

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$

Initialize action-value function $Q$ with random weights

**for** episode $= 1, M$ **do**

    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$

    **for** $t = 1, T$ **do**

        With probability $\epsilon$ select a random action $a_t$

        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$

        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$

        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

    **end for**

**end for**

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$

Initialize action-value function $Q$ with random weights

**for** episode $= 1, M$ **do**

    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced

    **for** $t = 1, T$ **do**

        With probability $\epsilon$ select a random action $a_t$

        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$

        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$

        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$

        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3

    **end for**

**end for**

## Deep Learning Framework Helps...

- Ideally, an end-user should only have
  - to give an environment, and
  - to choose a (deep) model

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q^*(\phi(s_t), a; \theta)$
        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$
        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$
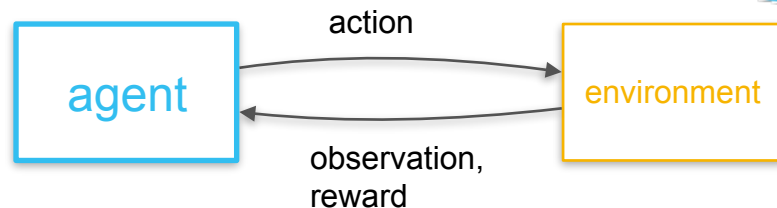        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$
        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
    **end for**
**end for**

Deep
**Reinforcement**
Learning
Framework Helps...

Preferred Networks

# Abstraction of training


action
observation, reward

- RL = <u>training loop</u>(environment, **agent**)
  - <u>training loop</u> alternately calls environment and agent

- **agent** = DQN(model, optimizer, explorer, replay_buf, phi)
  - In each step, agent returns an action
    after it updates the model
    - Agent.**act_and_train**(obs, reward) in ChainerRL

```
404|     def act_and_train(self, state, reward):
406|         with chainer.using_config('train', False):
407|             with chainer.no_backprop_mode():
408|                 action_value = self.model(self.batch_states([state], self.xp, self.phi))
411|                 greedy_action = cuda.to_cpu(action_value.greedy_actions.data)[0]

420|         action = self.explorer.select_action(self.t, lambda: greedy_action, action_value=action_value)
422|         self.t += 1

424|         # Update the target network
425|         if self.t % self.target_update_interval == 0:
426|             self.sync_target_network()

428|         if self.last_state is not None:
429|             assert self.last_action is not None
430|             # Add a transition to the replay buffer
431|             self.replay_buffer.append(
432|                 state=self.last_state, action=self.last_action, reward=reward,
435|                 next_state=state, next_action=action, is_state_terminal=False)
439|         self.last_state = state
440|         self.last_action = action

442|         self.replay_updater.update_if_necessary(self.t)

446|         return self.last_action
```

```
404|    def act_and_train(self, state, reward):
406|        with chainer.using_config('train', False):
407|            with chainer.no_backprop_mode():
408|                action_value = self.model(self.batch_states([state], self.xp, self.phi))
411|                greedy_action = cuda.to_cpu(action_value.greedy_actions.data)[0]
   |
420|        action = self.explorer.select_action(self.t, lambda: greedy_action, action_value=action_value)
422|        self.t += 1
   |
424|        # Update the target network
425|        if self.t % self.target_update_interval == 0:
426|            self.sync_target_network()
   |
428|        if self.last_state is not None:
429|            assert self.last_action is not None
430|            # Add a transition to the replay buffer
431|            self.replay_buffer.append(
432|                state=self.last_state, action=self.last_action, reward=reward,
435|                next_state=state, next_action=action, is_state_terminal=False)
439|        self.last_state = state
440|        self.last_action = action
   |
442|        self.replay_updater.update_if_necessary(self.t)
   |
446|        return self.last_action
```

exploration

Chainer RL    Preferred Networks

```
404|     def act_and_train(self, state, reward):
406|         with chainer.using_config('train', False):
407|             with chainer.no_backprop_mode():
408|                 action_value = self.model(self.batch_states([state], self.xp, self.phi))
411|                 greedy_action = cuda.to_cpu(action_value.greedy_actions.data)[0]
   |
420|         action = self.explorer.select_action(self.t, lambda: greedy_action,
422|         self.t += 1
   |
424|         # Update the target network
425|         if self.t % self.target_update_interval == 0:
426|             self.sync_target_network()
   |
428|         if self.last_state is not None:
429|             assert self.last_action is not None
430|             # Add a transition to the replay buffer
431|             self.replay_buffer.append(
432|                 state=self.last_state, action=self.last_action, reward=reward,
435|                 next_state=state, next_action=action, is_state_terminal=False)
439|         self.last_state = state
440|         self.last_action = action
   |
442|         self.replay_updater.update_if_necessary(self.t)
   |
446|         return self.last_action
```

replay memory

train

# Abstraction of Q-function & policy

- DQN model:
  - input: $s$ (state)
  - output: $[Q(s, a_1), \ldots, Q(s, a_n)]$ (Q-values)
- ActionValue object represents an output $Q(s, \text{-})$ of DQN
  - continuous action space $Q(s, a) = Q_0 - (a - a_0)^\top M (a - a_0)$ (NAF)

- Stochastic policy $[\pi(a_1 \mid s), \ldots, \pi(a_n \mid s)]$ (categorical distr.), $\pi(\text{-} \mid s) = N(\mu(s), \sigma^2(s))$, etc.

# Advanced techniques

- coroutine of exploring an episode
  - observation, reward = *yield* action
- logging Q, loss, etc.
  - chainer.report

# Sample code (dqn.ipynb)

# Exercise (exercise.ipynb)

- Double DQN [https://arxiv.org/abs/1509.06461v3]: eq. (4)

- Dueling Network [https://arxiv.org/abs/1511.06581v3]: eq. (9)

- Noisy Network [https://arxiv.org/abs/1706.10295v1]: replace linear layers (eq. (5)) with noisy linear layer (eq. (6))

  - **Hint:** Explorer := Greedy()

- Episodic replay (recurrent model, TD($\lambda$) target, etc.)

Gym environments with discrete actions (increasing order of difficulty)

- CartPole-v0 (return $\in$ [0, 200]), CartPole-v1 ([0, 500]), Acrobot-v1 ([-500, 0])

# References

- [Minh+ 2013] Mnih, Volodymyr, et al. "Playing atari with deep reinforcement learning." *arXiv preprint arXiv: 1312.5602* (2013).

  - Mnih, Volodymyr, et al. "Human-level control through deep reinforcement learning." *Nature* 518.7540 (2015): 529-533.

Preferred
Networks

**Algorithm 1** Deep Q-learning with Experience Replay

Initialize replay memory $\mathcal{D}$ to capacity $N$
Initialize action-value function $Q$ with random weights
**for** episode $= 1, M$ **do**
    Initialise sequence $s_1 = \{x_1\}$ and preprocessed sequenced $\phi_1 = \phi(s_1)$
    **for** $t = 1, T$ **do**
        With probability $\epsilon$ select a random action $a_t$
        otherwise select $a_t = \max_a Q(\phi(s_t), a; \theta)$
        Execute action $a_t$ in emulator and observe reward $r_t$ and image $x_{t+1}$
        Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$
        Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in $\mathcal{D}$
        Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from $\mathcal{D}$
        Set $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta') & \text{for non-terminal } \phi_{j+1} \end{cases}$

*$t_j$ is better*        $\theta'$ ("target network")

        Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ according to equation 3
    **end for**
**end for**

*this is usually called $y_j$*